

Dynamic Linking and Loading

\$Revision: 2.3 \$

\$Date: 1999/06/15 03:30:36 \$

Dynamic linking defers much of the linking process until a program starts running. It provides a variety of benefits that are hard to get otherwise:

- Dynamically linked shared libraries are easier to create than static linked shared libraries.
- Dynamically linked shared libraries are easier to update than static linked shared libraries.
- The semantics of dynamically linked shared libraries can be much closer to those of unshared libraries.
- Dynamic linking permits a program to load and unload routines at runtime, a facility that can otherwise be very difficult to provide.

There are a few disadvantages, of course. The runtime performance costs of dynamic linking are substantial compared to those of static linking, since a large part of the linking process has to be redone every time a program runs. Every dynamically linked symbol used in a program has to be looked up in a symbol table and resolved. (Windows DLLs mitigate this cost somewhat, as we describe below.) Dynamic libraries are also larger than static libraries, since the dynamic ones have to include symbol tables.

Beyond issues of call compatibility, a chronic source of problems is changes in library semantics. Since dynamic shared libraries are so easy to update compared to unshared or static shared libraries, it's easy to change libraries that are in use by existing programs, which means that the behavior of those programs changes even though "nothing has changed". This is a frequent source of problems on Microsoft Windows, where programs use a lot of shared libraries, libraries go through a lot of versions, and library version control is not very sophisticated. Most programs ship with copies of all of the libraries they use, and installers often will inadvertently install an older version of a shared library on top of a newer one, breaking programs that are expecting features found in the newer one. Well-behaved applications pop up a warning before installing an older library over a newer one, but even so, programs that depend on semantics of older libraries have been known to break when newer versions replace the older ones.

ELF dynamic linking

Sun Microsystems' SunOS introduced dynamic shared libraries to UNIX in the late 1980s. UNIX System V Release 4, which Sun co-developed, introduced the ELF object format and adapted the Sun scheme to ELF. ELF was clearly an improvement over the previous object formats, and by the late 1990s it had become the standard for UNIX and UNIX like systems including Linux and BSD derivatives.

Contents of an ELF file

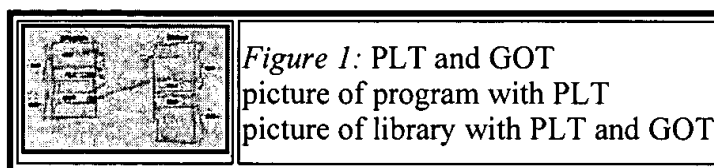
As mentioned in Chapter 3, an ELF file can be viewed as a set of *sections*, interpreted by the linker, or a set of *segments*, interpreted by the program loader. ELF programs and shared libraries have the same general structure, but with different sets of segments and sections.

ELF shared libraries can be loaded at any address, so they invariably use position independent code (PIC) so that the text pages of the file need not be relocated and can be shared among multiple processes. As described in Chapter 8, ELF linkers support PIC code with a Global Offset Table (GOT) in each shared library that contains pointers to all of the static data referenced in the program, Figure 1. The dynamic linker resolves and relocates all of the pointers in the GOT. This can be a performance issue but in practice the GOT is small except in very large libraries; a commonly

used version of the standard C library has only 180 entries in the GOT for over 350K of code.

Since the GOT is in the same loadable ELF file as the code that references it, and the relative addresses within a file don't change regardless of where the program is loaded, the code can locate the GOT with a relative address, load the address of the GOT into a register, and then load pointers from the GOT whenever it needs to address static data. A library need not have a GOT if it references no static data, but in practice all libraries do.

To support dynamic linking, each ELF shared library and each executable that uses shared libraries has a Procedure Linkage Table (PLT). The PLT adds a level of indirection for function calls analogous to that provided by the GOT for data. The PLT also permits "lazy evaluation", that is, not resolving procedure addresses until they're called for the first time. Since the PLT tends to have a lot more entries than the GOT (over 600 in the C library mentioned above), and most of the routines will never be called in any given program, that can both speed startup and save considerable time overall.



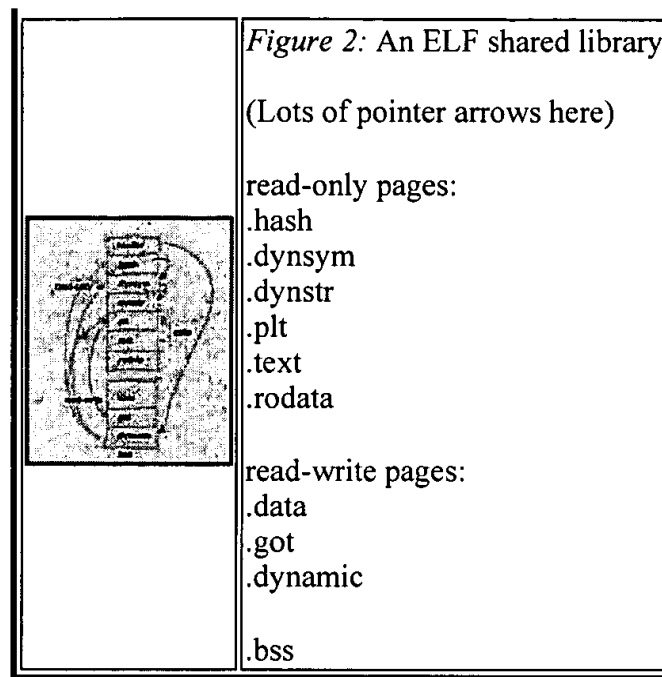
We discuss the details of the PLT below.

An ELF dynamically linked file contains all of the linker information that the runtime linker will need to relocate the file and resolve any undefined symbols. The `.dynsym` section, the dynamic symbol table, contains all of the file's imported and exported symbols. The `.dynstr` and `.hash` sections contain the name strings for the symbol, and a hash table the runtime linker can use to look up symbols quickly.

The final extra piece of an ELF dynamically linked file is the `DYNAMIC` segment (also marked as the `.dynamic` section) which runtime dynamic linker uses to find the information about the file the linker needs. It's loaded as part of the data segment, but is pointed to from the ELF file header so the runtime dynamic linker can find it. The `DYNAMIC` section is a list of tagged values and pointers. Some entry types occur just in programs, some just in libraries, some in both.

- **NEEDED**: the name of a library this file needs. (Always in programs, sometimes in libraries when one library is dependend on another, can occur more than once.)
- **SONAME**: "shared object name", the name of the file the linker uses. (Libraries.)
- **SYMTAB, STRTAB, HASH, SYMENT, STRSZ**: point to the symbol table, associated string and hash tables, size of a symbol table entry, size of string table. (Both.)
- **PLTGOT**: points to the GOT, or on some architectures to the PLT (Both.)
- **REL, RELSZ, and RELENT or RELA, RELASZ, and RELAENT**: pointer to, number of, and size of relocation entries. REL entries don't contain addends, RELA entries do. (Both.)
- **JMPREL, PLTRELSZ, and PLTREL**: pointer to, size, and format (REL or RELA) of relocation table for data referred to by the PLT. (Both.)
- **INIT and FINI**: pointer to initializer and finalizer routines to be called at program startup and finish. (Optional but usual in both.)
- A few other obscure types not often used.

An entire ELF shared library might look like Figure 2. First come the read-only parts, including the symbol table, PLT, text, and read-only data, then the read-write parts including regular data, GOT, and the dynamic section. The bss logically follows the last read-write section, but as always isn't present in the file.



An ELF program looks much the same, but in the read-only segment has `init` and `fini` routines, and an `INTERP` section near the front of the file to specify the name of the dynamic linker (usually `ld.so`). The data segment has no GOT, since program files aren't relocated at runtime.

Loading a dynamically linked program

Loading a dynamically linked ELF program is a lengthy but straightforward process.

Starting the dynamic linker

When the operating system runs the program, it maps in the file's pages as normal, but notes that there's an `INTERPRETER` section in the executable. The specified interpreter is the dynamic linker, `ld.so`, which is itself in ELF shared library format. Rather than starting the program, the system maps the dynamic linker into a convenient part of the address space as well and starts `ld.so`, passing on the stack an *auxiliary vector* of information needed by the linker. The vector includes:

- `AT_PHDR`, `AT_PHENT`, and `AT_PHNUM`: The address of the program header for the program file, the size of each entry in the header, and the number of entries. This structure describes the segments in the loaded file. If the system hasn't mapped the program into memory, there may instead be a `AT_EXECD` entry that contains the file descriptor on which the program file is open.
- `AT_ENTRY`: starting address of the program, to which the dynamic linker jumps after it has finished initialization.
- `AT_BASE`: The address at which the dynamic linker was loaded

At this point, bootstrap code at the beginning of `ld.so` finds its own GOT, the first entry in which points to the `DYNAMIC` segment in the `ld.so` file. From the dynamic segment, the linker can find its own relocation entries, relocate pointers in its own data segment, and resolve code references to the routines needed to load everything else. (The Linux `ld.so` names all of the essential routines with names starting with `_dt_` and special-case code looks for symbols that start with the string and resolves them.)

The linker then initializes a chain of symbol tables with pointers to the program's symbol table and the linker's own symbol table. Conceptually, the program file and all of the libraries loaded into a process share a single symbol table. But rather than build a merged symbol table at runtime, the linker keeps a linked list of the symbol tables in each file. Each file contains a hash table to speed symbol lookup, with a set of hash headers and a hash chain for each header. The linker can search for a symbol quickly by computing the symbol's hash value once, then running through appropriate hash chain in each of the symbol tables in the list.

Finding the libraries

Once the linker's own initializations are done, it finds the names of the libraries required by the program. The program's program header has a pointer to the "dynamic" segment in the file that contains dynamic linking information. That segment contains a pointer, `DT_STRTAB`, to the file's string table, and entries `DT_NEEDED` each of which contains the offset in the string table of the name of a required library.

For each library, the linker finds the library's ELF shared library file, which is in itself a fairly complex process. The library name in a `DT_NEEDED` entry is something like *libXt.so.6* (the Xt toolkit, version 6.) The library file might in in any of several library directories, and might not even have the same file name. On my system, the actual name of that library is `/usr/X11R6/lib/libXt.so.6.0`, with the ``.0` at the end being a minor version number.

The linker looks in these places to find the library:

- If the dynamic segment contains an entry called `DT_RPATH`, it's a colon-separated list of directories to search for libraries. This entry is added by a command line switch or environment variable to the regular (not dynamic) linker at the time a program is linked. It's mostly used for subsystems like databases that load a collection of programs and supporting libraries into a single directory.
- If there's an environment symbol `LD_LIBRARY_PATH`, it's treated as a colon-separated list of directories in which the linker looks for the library. This lets a developer build a new version of a library, put it in the `LD_LIBRARY_PATH` and use it with existing linked programs either to test the new library, or equally well to instrument the behavior of the program. (It skips this step if the program is set-uid, for security reasons.)
- The linker looks in the library cache file `/etc/ld.so.conf` which contains a list of library names and paths. If the library name is present, it uses the corresponding path. This is the usual way that most libraries are found. (The file name at the end of the path need not be exactly the same as the library name, see the section on library versions, below.)
- If all else fails, it looks in the default directory `/usr/lib`, and if the library's still not found, displays an error message and exits.

Once it's found the file containing the library, the dynamic linker opens the file, and reads the ELF header to find the program header which in turn points to the file's segments including the dynamic segment. The linker allocates space for the library's text and data segments and maps them in, along with zeroed pages for the bss. From the library's dynamic segment, it adds the library's symbol table to the chain of symbol tables, and if the library requires further libraries not already loaded, adds any new libraries to the list to be loaded.

When this process terminates, all of the libraries have been mapped in, and the loader has a logical global symbol table consisting of the union of all of the symbol tables of the program and the mapped library.

Shared library initialization

Now the loader revisits each library and handles the library's relocation entries, filling in the library's GOT and performing any relocations needed in the library's data segment. Load-time relocations on an x86 include:

- R_386_GLOB_DAT, used to initialize a GOT entry to the address of a symbol defined in another library.
- R_386_32, a non-GOT reference to a symbol defined in another library, generally a pointer in static data.
- R_386_RELATIVE, for relocatable data references, typically a pointer to a string or other locally defined static data.
- R_386_JMP_SLOT, used to initialize GOT entries for the PLT, described later.

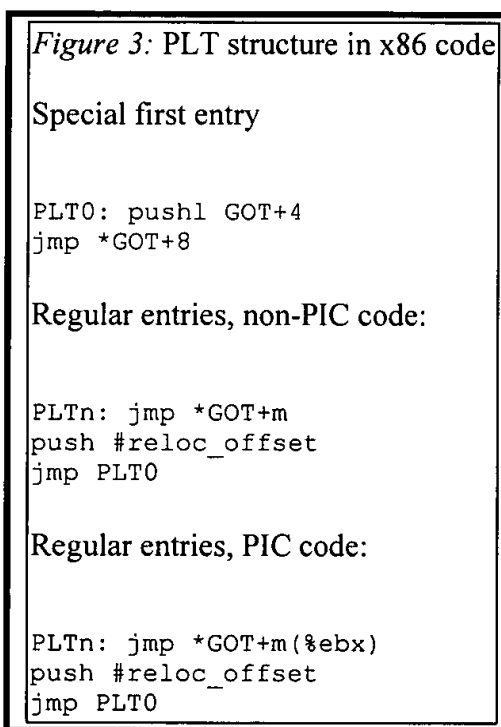
If a library has an `.init` section, the loader calls it to do library-specific initializations, such as C++ static constructors, and any `.fini` section is noted to be run at exit time. (It doesn't do the init for the main program, since that's handled in the program's own startup code.) When this pass is done, all of the libraries are fully loaded and ready to execute, and the loader calls the program's entry point to start the program.

Lazy procedure linkage with the PLT

Programs that use shared libraries generally contain calls to a lot of functions. In a single run of the program many of the functions are never called, in error routines or other parts of the program that aren't used. Furthermore, each shared library also contains calls to functions in other libraries, even fewer of which will be executed in a given program run since many of them are in routines that the program never calls either directly or indirectly.

To speed program startup, dynamically linked ELF programs use lazy binding of procedure addresses. That is, the address of a procedure isn't bound until the first time the procedure is called.

ELF supports lazy binding via the Procedure Linkage Table, or PLT. Each dynamically bound program and shared library has a PLT, with the PLT containing an entry for each non-local routine called from the program or library, Figure 3. Note that the PLT in PIC code is itself PIC, so it can be part of the read-only text segment.



All calls within the program or library to a particular routine are adjusted when the program or library is built to be calls to the routine's entry in the PLT. The first time the program or library calls a routine, the PLT entry calls the runtime linker to resolve the actual address of the routine. After that, the PLT entry jumps directly to the actual

address, so after the first call, the cost of using the PLT is a single extra indirect jump at a procedure call, and nothing at a return.

The first entry in the PLT, which we call PLT0, is special code to call the dynamic linker. At load time, the dynamic linker automatically places two values in the GOT. At GOT+4 (the second word of the GOT) it puts a code that identifies the particular library. At GOT+8, it puts the address of the dynamic linker's symbol resolution routine.

The rest of the entries in the PLT, which we call PLTn, each start with an indirect jump through a GOT entry. Each PLT entry has a corresponding GOT entry which is initially set to point to the push instruction in the PLT entry that follows the jmp. (In a PIC file this requires a loadtime relocation, but not an expensive symbol lookup.) Following the jump is a push instruction which pushes a relocation offset, the offset in the file's relocation table of a special relocation entry of type `R_386_JMP_SLOT`. The relocation entry's symbol reference points to the symbol in the file's symbol table, and its address points to the GOT entry.

This compact but rather baroque arrangement means that the first time the program or library calls a PLT entry, the first jump in the PLT entry in effect does nothing, since the GOT entry through which it jumps points back into the PLT entry. Then the push instruction pushes the offset value which indirectly identifies both the symbol to resolve and the GOT entry into which to resolve it, and jumps to PLT0. The instructions in PLT0 push another code that identifies which program or library it is, and then jumps into stub code in the dynamic linker with the two identifying codes at the top of the stack. Note that this was a jump, rather than a call, above the two identifying words just pushed is the return address back to the routine that called into the PLT.

Now the stub code saves all the registers and calls an internal routine in the dynamic linker to do the resolution. The two identifying words suffice to find the library's symbol table and the routine's entry in that symbol table. The dynamic linker looks up the symbol value using the concatenated runtime symbol table, and stores the routine's address into the GOT entry. Then the stub code restores the registers, pops the two words that the PLT pushed, and jumps off to the routine. The GOT entry having been updated, subsequent calls to that PLT entry jump directly to the routine itself without entering the dynamic linker.

Other peculiarities of dynamic linking

The ELF linker and dynamic linker have a lot of obscure code to handle special cases and try and keep the runtime semantics as similar as possible to those of unshared libraries.

Static initializations

If a program has an external reference to a global variable defined in a shared library, the linker has to create in the program a copy of the variable, since program data addresses have to be bound at link time, Figure 4. This poses no problem for the code in the shared library, since the code can refer to the variable via a GOT pointer which the dynamic linker can fix up, but there is a problem if the library initializes the variable. To deal with this problem, the linker puts an entry in the program's relocation table (which otherwise just contains `R_386_JMP_SLOT`, `R_386_GLOB_DAT`, `R_386_32`, and `R_386_RELATIVE` entries) of type `R_386_COPY` that points to the place in the program where the copy of the variable is defined, and tells the dynamic linker to copy the initial value of that word of data from the shared library.

Figure 4: Global data initialization

Main program:

```
extern int token;  
  
Routine in shared library:  
  
int token = 42;
```

Although this feature is essential for certain kinds of code, it occurs very rarely in practice. This is a band-aid, since it only works for single word data. The initializers that do occur are always pointers to procedures or other data, so the band-aid suffices.

Library versions

Dynamic libraries are generally named with major and minor versions numbers, like `libc.so.1.1` but programs should be bound only to major version numbers like `libc.so.1` since minor versions are supposed to be upward compatible.

To keep program loading reasonably fast, the system manager maintains a cache file containing the full pathname most recent version of each library, which is updated by a configuration program whenever a new library is installed.

To support this design, each dynamically linked library can have a "true name" called the *SONAME* assigned at library creation time. For example, the library called `libc.so.1.1` would have a SONAME of `libc.so.1`. (The SONAME defaults to the library's name.) When the linker builds a program that uses shared libraries, it lists the SONAMES of the libraries it used rather than the actual names of the libraries. The cache creation program scans all of the directories that contain shared libraries, finds all of the shared libraries, extracts the SONAME from each one, and where there are multiple libraries with the same SONAME, discards all but the highest version number. Then it writes the cache file with SONAMES and full pathnames so at runtime the dynamic linker can quickly find the current version of each library.

Dynamic loading at runtime

Although the ELF dynamic linker is usually called implicitly at program load time and from PLT entries, programs can also call it explicitly using `dlopen()` to load a shared library and `dlsym()` to find the address of a symbol, usually a procedure to call. Those two routines are actually simple wrappers that call back into the dynamic linker. When the dynamic linker loads a library via `dlopen()`, it does the same relocation and symbol resolution it does on any other library, so the dynamically loaded program can without any special arrangements call back to routines already loaded and refer to global data in the running program.

This permits users to add extra functionality to programs without access to the source code of the programs and without even having to stop and restart the programs (useful when the program is something like a database or a web server.) Mainframe operating systems have provided access to "exit routines" like this since at least the early 1960s, albeit without such a convenient interface, and it's long been a way to add great flexibility to packaged applications. It also provides a way for programs to extend themselves; there's no reason a program couldn't write a routine in C or C++, run the compiler and linker to create a shared library, then dynamically load and run the new code. (Mainframe sort programs have linked and loaded custom inner loop code for each sort job for decades.)

Microsoft Dynamic Link Libraries

Microsoft Windows also provides shared libraries, called dynamic-link libraries or DLLs in a fashion similar to but somewhat simpler than ELF shared libraries. The design of DLLs changed substantially between the 16 bit Windows 3.1 and the 32 bit Windows NT and 95. This discussion addresses only the more modern Win32 libraries. DLLs import procedure addresses using a PLT-like scheme. Although the design of DLLs would make it possible to import data addresses using a GOT-like scheme, in practice they use a simpler scheme that requires explicit program code to dereference imported pointers to shared data.

In Windows, both programs and DLLs are PE format (portable executable) files are intended to be memory mapped into a process. Unlike Windows 3.1, where all applications shared a single address space, Win32 gives each application its own address space and executables and libraries are mapped into each address space where they are used. For read-only code this doesn't make any practical difference, but for data it means that each application using a DLL gets its own copy of the DLL's data. (That's a slight oversimplification, since PE files can mark some sections as shared data with a single copy shared among all applications that use the file, but most data is unshared.)

Loading a Windows executable and DLLs is similar to loading a dynamically linked ELF program, although in the Windows case the dynamic linker is part of the kernel. First the kernel maps in the executable file, guided by section info in the PE headers. Then it maps in all of the DLLs that the executable refers to, again guided by the PE headers in each DLL.

PE files can contain relocation entries. An executable generally won't contain them and so has to be mapped at the address for which it was linked. DLLs all do contain relocation entries, and are relocated when they're mapped in if the address space for which they were linked isn't available. (Microsoft calls runtime relocation *rebasing*.)

All PE files, both executables and DLLs, have an entry point, and the loader calls a DLL's entry point when the DLL is loaded, when the DLL is unloaded, and each time a process thread attaches to or detaches from the DLL. (The loader passes an argument to say why it's making each call.) This provides a hook for static initializers and destructors analogous to the ELF `.init` and `.fini` sections.

Imported and exported symbols in PE files

PE supports shared libraries with two special sections of the file, `.edata`, for exported data, that lists the symbols exported from a file, and `.idata`, that lists the symbols imported into a file. Program files generally have only an `.idata` section, while DLLs always have an `.edata` and may have a `.idata` if they use other DLLs. Symbols can be exported either by symbol name, or by "ordinal", a small integer that gives the index of the symbol in the export address table. Linking by ordinals is slightly more efficient since it avoids a symbol lookup, but considerably more error prone since it's up to the person who builds a DLL to ensure that ordinals stay the same from one library version to another. In practice ordinals are usually used to call system services that rarely change, and names for everything else.

The `.edata` section contains an export directory table that describes the rest of the section, followed by the tables that define the exported symbols, Figure 5.

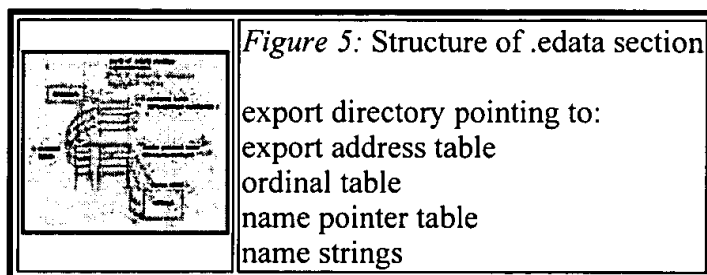
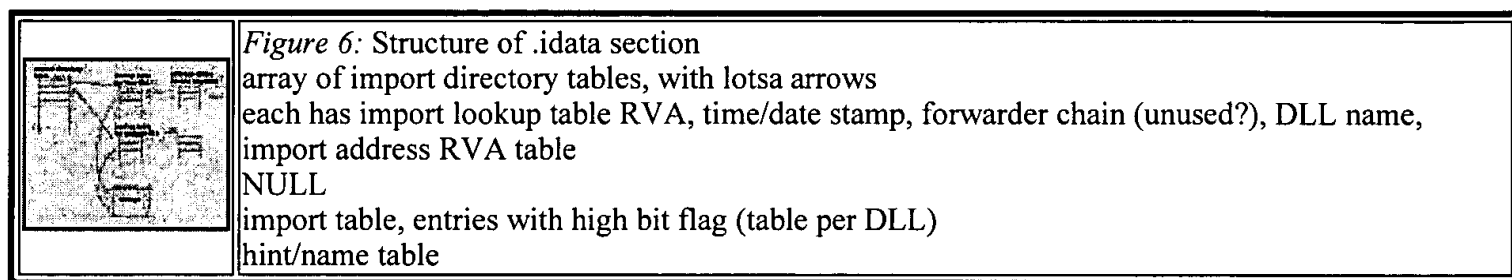


Figure 5: Structure of `.edata` section

The export address table contains the RVA (relative virtual address, relative to the base of the PE file) of the symbol. If the RVA points back into the `.edata` section, it's a "forwarder" reference, and the value pointed to is a string naming the symbol to use to satisfy the reference, probably defined in a different DLL. The ordinal and name pointer tables are parallel, with each entry in the name pointer table being the RVA of the name string for the symbol, and the ordinal being the index in the export address table. (Ordinals need not be zero-based; the ordinal base to subtract from ordinal values to get the index in the export address table is stored in the export directory and is most often 1.) Exported symbols need not all have names, although in practice they always do. The symbols in the name pointer table are in alphabetical order to permit the loader to use a binary search.

The `.idata` section does the converse of what the `.edata` section does, by mapping symbols or ordinals back into virtual addresses. The section consists of a null-terminated array of import directory tables, one per DLL from which symbols are imported, followed by an import lookup table per DLL, followed by a name table with hints, Figure 6.



For each imported DLL, there is an array of import addresses, typically in the program's text segment, into which the program loader places the resolved addresses. The import lookup table identifies the symbols to import, with the entries in the import lookup table being parallel to those in the import address table. The lookup table consists of 32 bit entries. If the high bit of an entry is set, the low 31 bits are the ordinal of the symbol to import, otherwise the entry is the RVA of an entry in the hint/name table. Each hint/name entry consists of a four-byte hint that guesses the index of the symbol in the DLL's export name pointer table, followed by the null terminated symbol name. The program loader uses the hint to probe the export table, and if the symbol name matches, it uses that symbol, otherwise it binary searches the entire export table for the name. (If the DLL hasn't changed, or at least its list of exported symbols hasn't changed, since the program that uses the DLL was linked, the guess will be right.)

Unlike ELF imported symbols, the values of symbols imported via `.idata` are only placed in the import address table, not fixed up anywhere else in the importing file. For code addresses, this makes little difference. When the linker builds an executable or DLL, it creates in the text section a table of misnamed "thunks", indirect jumps through the entries in the import address table, and uses the addresses of the thunks as the address of the imported routine, which is transparent to the programmer. (The thunks as well as most of the data in the `.idata` section actually come from a stub library created at the same time as the DLL.) In recent versions of Microsoft's C and C++ compiler, if the programmer knows that a routine will be called in a DLL, the routine can be declared "dllimport", and the compiler will emit an indirect call to the address table entry, avoiding the extra indirect jump. For data addresses, the situation is more problematical, since it's harder to hide the extra level of indirection required to address a symbol in another executable. Traditionally, programmers just bit the bullet and explicitly declared imported variables to be pointers to the real values and explicitly dereferenced the pointers. Recent versions of Microsoft's C and C++ compiler also let the programmer declare global data to be "dllimport" and the compiler will emit the extra pointer deferences, much like ELF code that references data indirectly via pointers in the GOT.

Lazy binding

Recent versions of Windows compilers have added delay loaded imports to permit lazy symbol binding for procedures, somewhat like the ELF PLT. A delay-loaded DLL has a structure similar to the `.idata` import directory

table, but not in the *.idata* section so the program loader doesn't handle it automatically. The entries in the import address table initially all point to a helper routine that finds and loads the DLL and replaces the contents of the address table with the actual addresses. The delay-loaded directory table has a place to store the original contents of the import address table so the values can be put back if the DLL is later unloaded. Microsoft provides a standard helper routine, but its interfaces are documented and programmers can write their own versions if need be.

Windows also permits programs to load and unload DLLs explicitly using `LoadLibrary` and `FreeLibrary`, and to find addresses of symbols using `GetProcAddress`.

DLLs and threads

One area in which the Windows DLL model doesn't work particularly well is thread local storage. A Windows program can start multiple threads in the same process, which share the process' address space. Each thread has a small chunk of thread local storage (TLS) to keep data specific to that thread, such as pointers to data structures and resources that the thread is using. The TLS needs "slots" for the data from the executable and from each DLL that uses TLS. The Windows linker can create a *.tls* section in a PE executable, that defines the layout for the TLS needed by routines in the executable and any DLLs to which it directly refers. Each time the process creates a thread, the new thread gets its own TLS, created using the *.tls* section as a template.

The problem is that most DLLs can either be linked implicitly from the executable, or loaded explicitly with `LoadLibrary`. DLLs loaded explicitly don't automatically get *.tls* storage, and since a DLL's author can't predict whether a library will be invoked implicitly or explicitly, it can't depend on the *.tls* section.

Windows defines runtime system calls that allocate slots at the end of the TLS. DLLs use those calls rather than *.tls* unless the DLL is known only to be invoked implicitly.

OSF/1 pseudo-static shared libraries

OSF/1, the ill-fated UNIX variant from the Open Software Foundation, used a shared library scheme intermediate between static and dynamic linking. Its authors noted that static linking is a lot faster than dynamic since less relocation is needed, and that libraries are updated infrequently enough that system managers are willing to endure some pain when they update shared libraries, although not the agony of relinking every executable program in the entire system.

So OSF/1 took the approach of maintaining a global symbol table visible to all processes, and loaded all the shared libraries into a sharable address space at system boot time. This assigned all of the libraries addresses that wouldn't change while the system was running. Each time a program started, if it used shared libraries, it would map in the shared libraries and symbol table and resolve undefined references in the executable using the global symbol table. No load-time relocation was ever required since programs were all linked to load in a part of the address space that was guaranteed to be available in each process, and the library relocation had already happened when they were loaded at boot time.

When one of the shared libraries changed, the system just had to be rebooted normally, at which point the system loaded the new libraries and created a new symbol table for executables to use.

This scheme was clever, but it wasn't very satisfactory. For one thing, processing symbol lookups is considerably slower than processing relocation entries, so avoiding relocation wasn't that much of a performance advantage. For another, dynamic linking provides the ability to load and run a library at runtime, and the OSF/1 scheme didn't provide for that.

Making shared libraries fast

Shared libraries, and ELF shared libraries in particular, can be very slow. The slowdowns come from a variety of sources, several of which we mentioned in Chapter 8:

- Load-time relocation of libraries
- Load-time symbol resolution in libraries and executables
- Overhead due to PIC function prolog code
- Overhead due to PIC indirect data references
- Slower code due to PIC reserved addressing registers

The first two problems can be ameliorated by caching, the latter two by retreating from pure PIC code.

On modern computers with large address spaces, it's usually possible to choose an address range for a shared library that's available in all or at least most of the processes that use the library. One very effective technique is similar to the Windows approach. Either when the library is linked or the first time a library is loaded, tentatively bind its addresses to a chunk of address space. After that, each time a program links to the library, use the same addresses of possible, which means that no relocation will be necessary. If that address space isn't available in a new process, the library is relocated as before.

SGI systems use the term *QUICKSTART* to describe the process of pre-relocating objects at linktime, or in a separate pass over the shared library. BeOS caches the relocated library the first time it's loaded into a process. If multiple libraries depend on each other, in principle it should be possible to pre-relocate and then pre-resolve symbol references among libraries, although I'm not aware of any linkers that do so.

If a system uses pre-relocated libraries, PIC becomes a lot less important. All the processes that load a library at its pre-relocated address can share the library's code whether it's PIC or not, so a non-PIC library at a well-chosen address can in practice be as sharable as PIC without the performance loss of PIC. This is basically the static linked library approach from Chapter 9, except that in case of address space collisions, rather than the program failing the dynamic linker moves the libraries at some loss of performance. Windows uses this approach.

BeOS implements cached relocated libraries with great thoroughness, including preserving correct semantics when libraries change. When a new version of a library is installed BeOS notes the fact and creates a new cached version rather than using the old cached version when programs refer to the library. Library changes can have a ripple effect. When library A refers to symbols in library B and B is updated, a new cached version of A will also have to be created if any of the referenced symbols in B have moved. This does make the programmer's life easier, but it's not clear to me that libraries are in practice updated often enough to merit the considerable amount of system code needed to track library updates.

Comparison of dynamic linking approaches

The Unix/ELF and Windows/PE dynamic linking differ in several interesting ways.

The ELF scheme uses a single name space per program, while the PE scheme uses a name space per library. An ELF executable lists the symbols it needs and the libraries it needs, but it doesn't record which symbol is in which library. A PE file, on the other hand, lists the symbols to import from each library. The PE scheme is less flexible but also more resistant to inadvertent spoofing. Imagine that an executable calls routine AFUNC which is found in library A and BFUNC which is found in library B. If a new version of library A happens to define its own BFUNC, an ELF program could use the new BFUNC in preference to the old one, while a PE program wouldn't. This is a problem with some large libraries; one partial solution is to use the poorly documented DT_FILTER and DT_AUXILIARY fields

to tell the dynamic linker what libraries this one imports symbols from, so the linker will search those libraries for imported symbols before searching the executable and the rest of the libraries. The `DT_SYMBOLIC` field tells the dynamic linker to search the library's own symbol table first, so that other libraries cannot shadow intra-library references. (This isn't always desirable; consider the malloc hack described in the previous chapter.) These ad-hoc approaches make it less likely that symbols in unrelated libraries will inadvertently shadow the correct symbols, but they're no substitute for a hierarchical link-time name space as we'll see in Chapter 11 that Java has.

The ELF scheme tries considerably harder than the PE scheme to maintain the semantics of static linked programs. In an ELF program, references to data imported from another library are automatically resolved, while a PE program needs to treat imported data specially. The PE scheme has trouble comparing the values of pointers to functions, since the address of an imported function is the address of the "thunk" that calls it, not the address of the actual function in the other library. ELF handles all pointers the same.

At run-time, nearly all of the Windows dynamic linker is in the operating system, while the ELF dynamic linker runs entirely as part of the application, with the kernel merely mapping in the initial files. The Windows scheme is arguably faster, since it doesn't have to map and relocate the dynamic linker in each process before it starts linking. The ELF scheme is definitely a lot more flexible. Since each executable names the "interpreter" program (now always the dynamic linker named `ld.so`) to use, different executables could use different interpreters without requiring any operating system changes. In practice, this makes it easier to support executables from variant versions of Unix, notably Linux and BSD, by making a dynamic linker that links to compatibility libraries that support non-native executables.

Exercises

In ELF shared libraries, libraries are often linked so that calls from one routine to another within a single shared library go through the PLT and have their addresses bound at runtime. Is this useful? Why or why not?

Imagine that a program calls a library routine `plugh()` that is found in a shared library, and the programmer builds a dynamically linked program that uses that library. Later, the system manager notices that `plugh` is a silly name for a routine and installs a new version of the library that calls the routine `xsazq` instead. What happens when the next time the programmer runs the program?

If the runtime environment variable `LD_BIND_NOW` is set, the ELF dynamic loader binds all of the program's PLT entries at load time. What would happen in the situation in the previous problem if `LD_BIND_NOW` were set?

Microsoft implemented lazy procedure binding without operating system assistance by adding some extra cleverness in the linker and using the existing facilities in the operating system. How hard would it be to provide transparent access to shared data, avoiding the extra level of pointers that the current scheme uses?

Project

It's impractical to build an entire dynamic linking system for our project linker, since much of the work of dynamic linking happens at runtime, not link time. Much of the work of building a shared library was already done in the project 8-3 that created PIC executables. A dynamically linked shared library is just a PIC executable with a well-defined list of imported and exported symbols and a list of other libraries on which it depends. To mark the file as a shared library or an executable that uses shared libraries, the first line is:

```
LINKLIB lib1 lib2 ...
or
```

```
LINK lib1 lib2 ...
```

where the lib's are the names of other shared libraries on which this one depends.

Project 10-1: Starting with the version of the linker from project 8-3, extend the linker to produce shared libraries and executables that need shared libraries. The linker needs to take as its input a list of input files to combine into the output executable or library, as well as other shared libraries to search. The output file contains a symbol table with defined (exported) and undefined (imported) symbols. Relocation types are the ones for PIC files along with AS4 and RS4 for references to imported symbols.

Project 10-2: Write a run-time binder, that is, a program that takes an executable that uses shared libraries and resolves its references. It should read in the executable, then read in the necessary libraries, relocating them to non-overlapping available addresses, and creating a logically merged symbol table. (You may want to actually create such a table, or use a list of per-file tables as ELF does.) Then resolve all of the relocations and external references. When you're done, all code and data should be assigned memory addresses, and all addresses in the code and data should be resolved and relocated to the assigned addresses.

[Next](#)
[Up](#)
[Previous](#)
[Contents](#)

Next: Unloading a Module Up: Modules Previous: Modules

Loading a Module

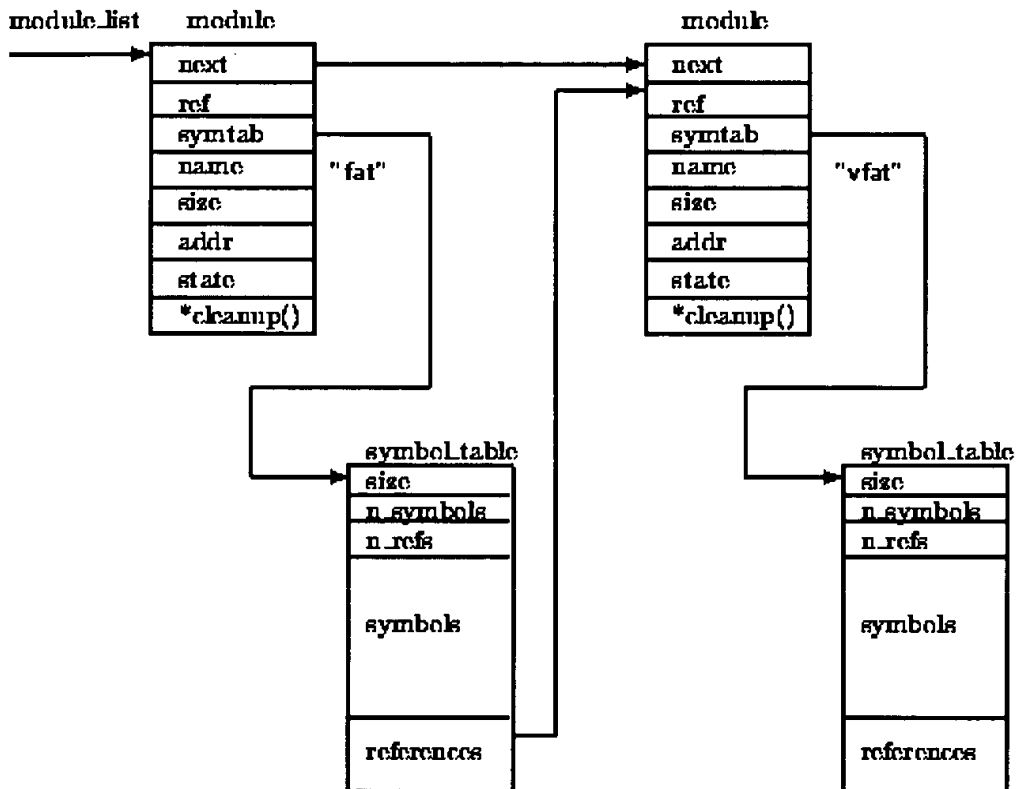


Figure: The List of Kernel Modules

There are two ways that a kernel module can be loaded. The first way is to use the *insmod* command to manually insert it into the kernel. The second, and much more clever way, is to load the module as it is needed; this is known as demand loading. When the kernel discovers the need for a module, for example when the user mounts a file system that is not in the kernel, the kernel will request that the kernel demon (*kernel*d) attempts to load the appropriate module.

The kernel demon is a normal user process albeit with super user privileges. When it is started up, usually at system boot time, it opens up an Inter-Process Communication (IPC) channel to the kernel. This link is used by the kernel to send messages to the *kernel*d asking for various tasks to be performed.

*Kernel*d's major function is to load and unload kernel modules but it is also capable of other tasks such as starting up the PPP link over serial line when it is needed and closing it down when it is not. *Kernel*d does not perform these tasks itself, it runs the necessary programs such as *insmod* to do the work. *Kernel*d is just an agent of the kernel, scheduling work on its behalf.


The *insmod* utility must find the requested kernel module that it is to load. Demand loaded kernel modules are normally kept in `/lib/modules/kernel-version`. The kernel modules are linked object files just like other

programs in the system except that it is linked as a relocatable image. That is an image that is not linked to run from a particular address. They can be either `a.out` or `elf` format object files. *Insmod* makes a privileged system call to find the kernel's exported symbols.


These are kept in pairs containing the symbol's name and its value, for example its address. The kernel's exported symbol table is held in the first `module` data structure in the list of modules maintained by the kernel and pointed at by the `module_list` pointer.

Only specifically entered symbols are added into the table which is built when the kernel is compiled and linked, not every symbol in the kernel is exported to its modules. An example symbol is `request_irq` which is the kernel routine that must be called when a driver wishes to take control of a particular system interrupt. In my current kernel, this has a value of `0x0010cd30`. You can easily see the exported kernel symbols and their values by looking at `/proc/ksyms` or by using the *ksyms* utility. The *ksyms* utility can either show you all of the exported kernel symbols or only those symbols exported by loaded modules. *Insmod* reads the module into its virtual memory and fixes up its unresolved references to kernel routines and resources using the exported symbols from the kernel. This fixing up takes the form of patching the module image in memory. *Insmod* physically writes the address of the symbol into the appropriate place in the module.

When *insmod* has fixed up the module's references to exported kernel symbols, it asks the kernel for enough space to hold the new kernel, again using a privileged system call. The kernel allocates a new `module` data structure and enough kernel memory to hold the new module and puts it at the end of the kernel modules list. The new module is marked as `UNINITIALISED`.

Figure  shows the list of kernel modules after two modules, `VFAT` and `VFAT` have been loaded into the kernel. Not shown in the diagram is the first module on the list which is a pseudo-module that is only there to hold the kernel's exported symbol table. You can use the command *lsmod* to list all of the loaded kernel modules and their interdependencies. *Lsmod* simply reformats `/proc/modules` which is built from the list of kernel `module` data structures. The memory that the kernel allocates for it is mapped into the *insmod* process's address space so that it can access it. *Insmod* copies the module into the allocated space and relocates it so that it will run from the kernel address that it has been allocated. This must happen as the module cannot expect to be loaded at the same address twice let alone into the same address in two different Linux systems. Again, this relocation involves patching the module image with the appropriate addresses.

The new module also exports symbols to the kernel and *insmod* builds a table of these exported images. Every kernel module must contain module initialization and module cleanup routines and these symbols are deliberately not exported but *insmod* must know the addresses of them so that it can pass them to the kernel. All being well, *insmod* is now ready to initialize the module and it makes a privileged system call passing the kernel the addresses of the module's initialization and cleanup routines.

When a new module is added into the kernel, it must update the kernel's set of symbols and modify the modules that are being used by the new module. Modules which have other modules dependent on them must maintain a list of references at the end of their symbol table and pointed at by their `module` data structure. Figure  shows that the `VFAT` file system module is dependent on the `FAT` file system module. So, the `FAT` module contains a reference to the `VFAT` module; the reference was added when the `VFAT` module was loaded. The kernel calls the modules initialization routine and, if it is successful it carries on installing the module. The module's cleanup routine address is stored in its `module` data structure and it will be called by the kernel when that module is unloaded. Finally, the module's state is set to `RUNNING`.

[Next](#) [Up](#) [Previous](#) [Contents](#)

Next: Unloading a Module **Up:** Modules **Previous:** Modules

David A. Rusling

david.rusling@reo.mts.dec.com

[Next](#) [Up](#) [Previous](#) [Contents](#)

Next: The Linux Kernel Sources **Up:** Modules **Previous:** Loading a Module

Unloading a Module

Modules can be removed using the *rmmod* command but demand loaded modules are automatically removed from the system by *kerneld* when they are no longer being used. Every time its idle timer expires, *kerneld* makes a system call requesting that all unused demand loaded modules are removed from the system. The timer's value is set when you start *kerneld*, my *kerneld* checks every 180 seconds. So, for example, if you mount an *iso9660* CD ROM and your *iso9660* filesystem is a loadable module, then shortly after the CD ROM is unmounted, the *iso9660* module will be removed from the kernel.

A module cannot be unloaded so long as other components of the kernel are depending on it. For example, you cannot unload the *VFAT* module if you have one or more *VFAT* file systems mounted. If you look at the output of *lsmod*, you will see that each module has a count associated with it. For example:

```
Module:           #pages:  Used by:
msdos              5              1
vfat               4              1 (autoclean)
fat               6      [vfat msdos] 2 (autoclean)
```

The count is the number of kernel entities that are dependent on this module. In the above example, the *vfat* and *msdos* modules are both dependent on the *fat* module and so it has a count of 1. Both the *vfat* and *msdos* modules have 1 dependent, which is a mounted file system. If I were to load another *VFAT* file system then the *vfat* module's count would become 2. A module's count is held in the first longword of its image.

This field is slightly overloaded as it also holds the *AUTOCLEAN* and *VISITED* flags. Both of these flags are used for demand loaded modules. These modules are marked as *AUTOCLEAN* so that the system can recognize which ones it may automatically unload. The *VISITED* flag marks the module as in use by one or more other system components; it is set whenever another component makes use of the module. Each time the system is asked by *kerneld* to remove unused demand loaded modules it looks through all of the modules in the system for likely candidates. It only looks at modules marked as *AUTOCLEAN* and in the state *RUNNING*. If the candidate has its *VISITED* flag cleared then it will remove the module, otherwise it will clear the *VISITED* flag and go on to look at the next module in the system.

Assuming that a module can be unloaded, its cleanup routine is called to allow it to free up the kernel resources that it has allocated.

The *module* data structure is marked as *DELETED* and it is unlinked from the list of kernel modules. Any other modules that it is dependent on have their reference lists modified so that they no longer have it as a dependent. All of the kernel memory that the module needed is deallocated.

[Next](#) [Up](#) [Previous](#) [Contents](#)

Next: The Linux Kernel Sources **Up:** Modules **Previous:** Loading a Module

David A. Rusling
david.rusling@reo.mts.dec.com

Introduction to Linux kernel modules

By Vans Information <Content@vansinfo.com>

Posted: (2001-05-16 07:56:03 EST by gatha)

The earlier articles touched on what Linux Operating System is all about. You know its scope, its advantages, disadvantages, weaknesses and strongholds. In this article, we take a closer look at the Kernel modules--one of the two modes that Linux operates in.

Linux operates in two modes--the *Kernel mode* (kernel space) and the *User mode* (user space). The kernel works in the highest level (also called *supervisor mode*) where it has all the authority, while the applications work in the lowest level where direct access to hardware and memory are prohibited. Keeping in line with the traditional Unix philosophy, Linux transfers the execution from user space to the kernel space through system calls and the hardware interrupts. The Kernel code executing the system call works in the context of the process, which invokes the system call. As it operates on behalf of the calling process, it can access the data in the processes address space. The kernel code that handles interrupts, works to the processes and related to any particular process.

Linux Kernel Modules

The Linux kernel is a *monolithic kernel* i.e. it is one single large program where all the functional components of the kernel have access to all of its internal data structures and routines. The alternative to this is the *micro kernel* structure where the functional pieces of the kernel are broken out into units with strict communication mechanism between them. This makes adding new components into the kernel, via the configuration process, rather time consuming. The best and the robust alternative is the ability to dynamically load and unload the components of the operating system using *Linux Kernel Modules*.

The Linux kernel modules are piece of codes, which can be dynamically linked to the kernel (according to the need), even after the system bootup. They can be unlinked from the kernel and removed when they are no longer needed. Mostly the Linux kernel modules are used for *device drivers* or *pseudo-device drivers* such as network drivers or file system. When a Linux kernel module is loaded, it becomes a part of the Linux kernel as the normal kernel code and functionality and it posses the same rights and responsibilities as the kernel code.

Life cycle of Linux kernel module

The life cycle of a module starts with the *init_module()*. The task of *init_module* is to prepare the module for later invocation. The module is registered to the kernel and attaches its data-structures and functionality to the kernel. The kernel-defined external functions are also resolved. The life cycle of the module ends with *cleanup_module()*. It 'unregisters' the module functionality from the kernel. We shall later take a more detailed look at the life cycle of the module.

Simple Module Program

Let us now program a simple module to review its life cycle. The *init_module* is called when the module is inserted into the kernel where as the *cleanup_module* is called just before removing it from the kernel. In the following program, the *init_module* and the *cleanup_module* functions are demonstrated.

```
/* Simple Linux kernel module Feb'2001
```

```
#include
#include
```

```

#if CONFIG_MODVERSIONS==1
#define MODVERSINS
#include

#endif

/ initialise the module /
int init_module()
{
  printk("init_module invoked\n");
  printk("the message is printed from the kernel space\n");
  / if the non zero value is returned, then it means that the init_module failed and the kernel module can't be loaded /

  return 0;
}

/ cleanup / end of module life cycle */
void cleanup_module()
{
  printk("cleanup_module invoked\n");
  printk("module is now going to be unloaded from the kernel\n");
}

```

Compile the above program using the following:

```
#gcc -Wall -DMODULE -D__KERNEL__ -DLINUX -O -c simpleModule.c
```

Run the compiled module using the following:

```
#insmod simpleModule.o
```

Remember, you have to run the above command from the Linux shell at raw console (not from the console in Xwindows environment) at root login.

Now check the status of the module using

```
#lsmod
```

Then remove the module using **#rmmod simpleModule**

If you have not seen any of the module-initiated console printing (implemented using `printk') about the status of the module, use the following command to see the kernel messages. **dmesg | less** In the above, commands **insmod**, **lsmod** and **rmmod** are used to load and unload modules to the Linux kernel. The details are discussed in the following section.

Loading modules **insmod** loads the 'loadable kernel modules' in the running kernel. **insmod** tries to link a module into the running kernel by resolving all the symbols from the kernel's exported 'symbol table'.

Lets now discuss the demand loading of the module by the kernel, *dynamically*. When the Linux kernel discovers the need for a module, the kernel requests to the kernel daemon (kernelld) to load the appropriate module. To illustrate this with an example, lets mount a NTFS partition in the Linux system. If the NTFS filesystem support is not

statically implemented in the kernel (but compiled as a module), the kernel daemon will search for the appropriate module and load it from the repository. Then the partition is mounted for the use.

Lets go deep into the action of the kernel daemon (kerneld). The kerneld is the normal user process having exclusive superuser privileges. At the time of booting, kerneld opens the IPC channel to the kernel and uses it for transferring messages (request for loading modules), to and from the kernel. While loading the module, the kerneld calls *modprobe* and *insmod* to load the required module. The *insmod* utility should be able to access the requested module. The demand loadable kernel modules are usually located at */lib/module/* directory as the *object files linked as relocatable images*.

Let us revisit the working of the *insmod* to get a clear picture of the module loading operation. The *insmod* depends on some critical system calls to load the module to the kernel. It uses the *sys_create_module* to allocate kernel memory to hold module. It uses *get_kernel_syms* system call to get the kernel symbol table in order to link the module. It then calls the *sys_init_module* system call to copy the relocatable object code of the module to the kernel space. And soon after this, *insmod* calls the initialization function of the concerned module i.e. *init_module*. All of these system calls can be found in *kernel/module.c*.

Unloading modules

The modules can be unloaded using *rmmod* command. While removing modules, *rmmod* ensures the restriction that the modules are not in use and they are not referred by any other module or the part of the kernel. The *demand loaded modules* (i.e. the modules loaded by kerneld) are automatically removed from the system by 'kerneld' when they are no longer used. Every time it's idle, timer expires and kerneld makes a system call requesting for all the *demand loaded kernels*, which are not busy to be removed. The modules, whose visited flags are cleared and marked as *AUTOCLEAN*, are 'unloaded'.

Assuming that the module can be unloaded, the *cleanup_module* function of the concerned module is called to free up the kernel resources it has allocated. After the successful execution of the *cleanup_module*, the module datastructure is marked *DELETED* and it is unlinked from the kernel and unlisted from the list of kernel modules. The reference list of the modules on which it (module removed) is dependent is modified and dependency is released. The kernel memory allocated to the concerned module is deallocated and returned to the kernel memory pool.

Version Dependency of modules Version dependency of the module is one of the trickiest parts of the Linux Kernel Module programming. Typically, the modules are required to be compiled for each version of the kernel. Each module defines a symbol called *kernel_version*, which *insmod* matches against the version number of the current kernel. The current kernel 2.2.x/2.4.x define the symbol in . Hence if the module is made up of multiple source files, the should be included in only one of the source files.

Though typically, modules should be recompiled for each kernel version, it is not always possible to recompile module when it is run on as a commercial module distributed in binary form. Kernel developers have provided a flexible way to deal with the version problem. The idea is that a module is incompatible with a different kernel version only if the software interface offered by the kernel is changed. The software interface is represented by the function prototype and the exact definition of all the data structures involved in the function call. The CRC algorithm can be used to map all the information about software interface to the single 32bit number. The issue of version dependency is handled by using the name of the each symbol exported.

Other articles by Vans Information

© 1998-2004 FreeOS Technologies (I) Pvt. Ltd. All rights reserved. [Privacy Policy]

WebHosting Provided By Interpole Solutions

